

Anderson acceleration of coordinate descent

Quentin Bertrand¹

Mathurin Massias²

¹Université Paris-Saclay, Inria, CEA, Palaiseau, France

²MaLGA, DIBRIS, Università degli Studi di Genova

Abstract

Acceleration of first order methods is mainly obtained via inertial techniques à la Nesterov, or via nonlinear extrapolation. The latter has known a recent surge of interest, with successful applications to gradient and proximal gradient techniques. On multiple Machine Learning problems, coordinate descent achieves performance significantly superior to full-gradient methods. Speeding up coordinate descent in practice is not easy: inertially accelerated versions of coordinate descent are theoretically accelerated, but might not always lead to practical speed-ups. We propose an accelerated version of coordinate descent using extrapolation, showing considerable speed up in practice, compared to inertial accelerated coordinate descent and extrapolated (proximal) gradient descent. Experiments on least squares, Lasso, elastic net and logistic regression validate the approach.

1 Introduction

Gradient descent is the workhorse of modern convex optimization (Nesterov, 2004; Beck, 2017). For composite problems, proximal gradient descent retains the nice properties enjoyed by the latter. In both techniques, inertial acceleration achieves optimal convergence rates (Nesterov, 1983; Beck and Teboulle, 2009).

Coordinate descent is a variant of gradient descent, which updates the iterates one coordinate at a time (Tseng and Yun, 2009; Friedman et al., 2010). Proximal coordinate descent has been applied to numerous Machine Learning problems (Shalev-Shwartz and Zhang, 2013; Wright, 2015; Shi et al., 2016), in particular the Lasso (Tibshirani, 1996), elastic net (Zou and Hastie, 2005) or sparse logistic regression (Ng, 2004). It is used in preminent packages such as scikit-learn (Pedregosa et al., 2011), glmnet (Friedman et al.,

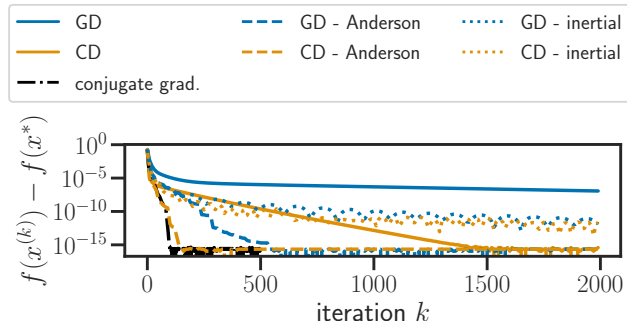


Figure 1: Suboptimality along iterations on the 1000 first features of the *rcv1* dataset for a quadratic problem. GD: gradient descent, CD: coordinate descent.

2009), libsvm (Fan et al., 2008) or lightning (Blondel and Pedregosa, 2016). On the theoretical side, inertial accelerated versions of coordinate descent (Nesterov, 2012; Lin et al., 2014; Fercoq and Richtárik, 2015) achieves optimal rates.

To obtain optimal rates, Anderson acceleration (Anderson, 1965) is an alternative to inertia: it provides acceleration by exploiting the iterates' structure. This procedure has been known for a long time, under various names and variants (Wynn, 1962; Eddy, 1979; Smith et al., 1987), see Sidi (2017); Brezinski et al. (2018) for reviews. Anderson acceleration enjoys accelerated rates on quadratic functions (Golub and Varga, 1961), but theoretical guarantees in the nonquadratic case are weaker (Scieur et al., 2016). Interestingly, numerical performances still show significant improvements on nonquadratic objectives. Anderson acceleration has been adapted to various algorithms such as Douglas-Rachford (Fu et al., 2019), ADMM (Poon and Liang, 2019) or proximal gradient descent (Zhang et al., 2018; Mai and Johansson, 2019; Poon and Liang, 2020). Among main benefits, the practical version of

Algorithm 1 Offline Anderson extrapolation

```

init:  $x^{(0)} \in \mathbb{R}^p$ 
1 for  $k = 1, \dots$  do
2    $x^{(k)} = Tx^{(k-1)} + b$  // regular linear iteration
3    $U = [x^{(1)} - x^{(0)}, \dots, x^{(k)} - x^{(k-1)}]$ 
4    $c = (U^\top U)^{-1} \mathbf{1}_k / \mathbf{1}_k^\top (U^\top U)^{-1} \mathbf{1}_k \in \mathbb{R}^k$ 
5    $x_{\text{e-off}}^{(k)} = \sum_{i=1}^k c_i x^{(i)}$  // does not affect  $x^{(k)}$ 
6 return  $x_{\text{e-off}}^{(k)}$ 

```

Anderson acceleration is memory efficient, easy to implement, line search free, has a low cost per iteration and does not require knowledge of the strong convexity constant. Finally, it introduces a single additional parameter, which often does not require tuning (see Section 3.1).

In this work:

- We propose an Anderson acceleration scheme for coordinate descent, which, as visible on Figure 1, outperforms inertial and extrapolated gradient descent, as well as inertial coordinate descent, even reaching the first order optimal performance of conjugate gradient on this example.
- The acceleration is obtained eventhough the iteration matrix is not symmetric, a notable problem in the analysis of Anderson extrapolation.
- We empirically highlight that the proposed acceleration technique can generalize in the non-quadratic case (Algorithm 3) and significantly improve proximal coordinate descent algorithms (Section 3), which are state-of-the-art first order methods on the considered problems.

Notation The j -th line of the matrix A is A_j ; and its j -th column is $A_{\cdot j}$. The canonical basis vectors of \mathbb{R}^p are e_j . The vector of size K with all one entries is $\mathbf{1}_K$. The spectral radius of the matrix A , $\rho(A)$, is the largest eigenvalue modulus of A . The set of p by p symmetric semidefinite positive matrices is \mathbb{S}_+^p . The condition number $\kappa(A)$ of a matrix A is its largest singular value divided by its smallest. A positive definite matrix A induces the norm $\|x\|_A = \sqrt{x^\top A x}$. The proximity operator of the function g is $\text{prox}_g(x) = \arg \min_y g(y) + \frac{1}{2} \|x - y\|^2$.

2 Anderson extrapolation

2.1 Background

Anderson extrapolation is designed to accelerate the convergence of sequences based on fixed point linear

Algorithm 2 Online Anderson extrapolation

```

init:  $x^{(0)} \in \mathbb{R}^p$ 
1 for  $k = 1, \dots$  do
2    $x^{(k)} = Tx^{(k-1)} + b$  // regular iteration
3   if  $k = 0 \pmod K$  then
4      $U = [x^{(k-K+1)} - x^{(k-K)}, \dots, x^{(k)} - x^{(k-1)}]$ 
5      $c = (U^\top U)^{-1} \mathbf{1}_K / \mathbf{1}_K^\top (U^\top U)^{-1} \mathbf{1}_K \in \mathbb{R}^K$ 
6      $x_{\text{e-on}}^{(k)} = \sum_{i=1}^K c_i x^{(k-K+i)}$ 
7      $x^{(k)} = x_{\text{e-on}}^{(k)}$  // base sequence changes
8 return  $x^{(k)}$ 

```

iterations, that is:

$$x^{(k+1)} = Tx^{(k)} + b, \quad (1)$$

where the *iteration matrix* $T \in \mathbb{R}^{p \times p}$ has spectral radius $\rho(T) < 1$. There exist two variants: offline and online, which we recall briefly.

Offline extrapolation (Algorithm 1), at iteration k , looks for a fixed point as an affine combination of the k first iterates: $x_{\text{e-off}}^{(k)} = \sum_{i=1}^k c_i x^{(i-1)}$, and solves for the coefficients $c^{(k)} \in \mathbb{R}^k$ as follows:

$$\begin{aligned}
c^{(k)} &= \arg \min_{\sum_{i=1}^k c_i = 1} \left\| \sum_{i=1}^k c_i x^{(i-1)} - T \sum_{i=1}^k c_i x^{(i-1)} - b \right\|^2 \\
&= \arg \min_{\sum_{i=1}^k c_i = 1} \left\| \sum_{i=1}^k c_i (x^{(i)} - x^{(i-1)}) \right\|^2 \\
&= (U^\top U)^{-1} \mathbf{1}_k / \mathbf{1}_k^\top (U^\top U)^{-1} \mathbf{1}_k, \quad (2)
\end{aligned}$$

where $U = [x^{(1)} - x^{(0)}, \dots, x^{(k)} - x^{(k-1)}] \in \mathbb{R}^{p \times k}$ (and hence the objective rewrites $\|Uc\|^2$). In practice, since $x^{(k)}$ is available when $c^{(k)}$ is computed, one uses $x_{\text{e}}^{(k)} = \sum_{i=1}^k c_i x^{(i)}$ instead of $\sum_{i=1}^k c_i x^{(i-1)}$. The motivation for introducing the coefficients $c^{(k)}$ is discussed in more depth after Prop. 6 in Massias et al. (2019), and details about the closed-form solution can be found in Scieur et al. (2016, Lem. 2.4). In offline acceleration, more and more base iterates are used to produce the extrapolated point, but the extrapolation sequence does not affect the base sequence. This may not scale well since it requires solving larger and larger linear systems.

A more practical variant is the *online* version (Algorithm 2), considered in this paper. The number of points to be extrapolated is fixed to K ; $x^{(1)}, \dots, x^{(K)}$ are computed normally with the fixed point iterations, but $x_{\text{e}}^{(K)}$ is computed by extrapolating the iterates from $x^{(1)}$ to $x^{(K)}$, and $x^{(K)}$ is taken equal to $x_{\text{e}}^{(K)}$. K normal iterates are then computed from $x^{(K+1)}$ to $x^{(2K)}$ then extrapolation is performed on these last K iterates, etc.

As we recall below, results on Anderson acceleration mainly concern fixed-point iterations with symmetric iteration matrices T , and results concerning

non-symmetric iteration matrices are weaker (Bollapragada et al., 2018). Poon and Liang (2020, Thm 6.4) do not assume that T is symmetric, but only diagonalizable, which is still a strong requirement.

Proposition 1 (Symmetric T , Scieur 2019). Let the iteration matrix T be symmetric semi-definite positive, with spectral radius $\rho = \rho(T) < 1$. Let x^* be the limit of the sequence $(x^{(k)})$. Let $\zeta = (1 - \sqrt{1 - \rho}) / (1 + \sqrt{1 - \rho})$. Then the iterates of offline Anderson acceleration satisfy, with $B = (\text{Id} - T)^2$:

$$\|x_{\text{e-off}}^{(k)} - x^*\|_B \leq \frac{2\zeta^{k-1}}{1 + \zeta^{2(k-1)}} \|x^{(0)} - x^*\|_B, \quad (3)$$

and thus those of online extrapolation satisfy:

$$\|x_{\text{e-on}}^{(k)} - x^*\|_B \leq \left(\frac{2\zeta^{K-1}}{1 + \zeta^{2(K-1)}} \right)^{k/K} \|x^{(0)} - x^*\|_B. \quad (4)$$

Scieur et al. (2016) showed that the offline version in Proposition 1 matches the accelerated rate of the conjugate gradient (Hestenes and Stiefel, 1952). As it states, gradient descent can be accelerated by Anderson extrapolation on quadratics.

Application to least squares The canonical application of Anderson extrapolation is gradient descent on least squares. Consider a quadratic problem, with $b \in \mathbb{R}^p$, $H \in \mathbb{S}_{++}^p$ such that $0 \prec H \preceq L$ and $L > 0$:

$$x^* = \arg \min_{x \in \mathbb{R}^p} \frac{1}{2} x^\top H x + \langle b, x \rangle. \quad (5)$$

A typical instance is overdetermined least squares with full-column rank design matrix $A \in \mathbb{R}^{n \times p}$, and observations $y \in \mathbb{R}^n$, such that $H = A^\top A$ and $b = -A^\top y$. On Problem (5) gradient descent with step size $1/L$ reads:

$$x^{(k+1)} = \underbrace{\left(\text{Id}_p - \frac{1}{L} H \right)}_{T^{\text{GD}} \in \mathbb{S}_+^p} x^{(k)} + \underbrace{(-b/L)}_{b^{\text{GD}}}. \quad (6)$$

Because they have this linear structure, iterates of gradient descent can benefit from Anderson acceleration, observing that the fixed point of $x \mapsto T^{\text{GD}} x + b^{\text{GD}}$ solves (5), with $T^{\text{GD}} \in \mathbb{S}_+^p$. Anderson acceleration of gradient descent has therefore been well-studied beyond the scope of Machine Learning (Pulay, 1980; Eyert, 1996). However, on many Machine Learning problems, coordinate descent achieves far superior performance, and it is interesting to determine whether or not it can also benefit from Anderson extrapolation.

2.2 Linear iterations of coordinate descent

To apply Anderson acceleration to coordinate descent, we need to show that its iterates satisfy linear iterations as in (6). An epoch of cyclic coordinate descent

for Problem (5) consists in updating the vector x one coordinate at a time, sequentially, i.e. for $j = 1, \dots, p$:

$$x_j \leftarrow x_j - \frac{1}{H_{jj}} (H_{j\cdot} x + b_j), \quad (7)$$

which can be rewritten, for $j = 1, \dots, p$:

$$x \leftarrow \left(\text{Id}_p - \frac{e_j e_j^\top}{H_{jj}} H \right) x - \frac{b_j}{H_{jj}} e_j. \quad (8)$$

Thus, for primal iterates, as observed by Bertrand et al. (2020, Sec. A.3), one full pass (updating coordinates from 1 to p) leads to a linear iteration:

$$x^{(k+1)} = T^{\text{CD}} x^{(k)} + b^{\text{CD}}, \quad (9)$$

with $T^{\text{CD}} = \left(\text{Id}_p - \frac{e_p e_p^\top}{H_{pp}} H \right) \dots \left(\text{Id}_p - \frac{e_1 e_1^\top}{H_{11}} H \right)$. Note that in the case of coordinate descent we write $x^{(k)}$ for the iterates after one pass of coordinate descent on all features, and not after each update (7). The iterates of coordinate descent therefore also have a fixed-point structure, but contrary to gradient descent, their iteration matrix T^{CD} is not symmetric, which we address in Section 2.3.

2.3 Anderson extrapolation for nonsymmetric iteration matrices

Even on quadratics, Anderson acceleration with non-symmetric iteration matrices is less developed, and the only results concerning its theoretical acceleration are recent and weaker than in the symmetric case.

Proposition 2 (Bollapragada et al. 2018, Thm 2.2). When T is not symmetric, and $\rho(T) < 1$,

$$\|x_{\text{e-off}}^{(k)} - T x_{\text{e-off}}^{(k)} - b\| \leq \|\text{Id} - \rho(T - \text{Id})\|_2 \|P^*(T)(x^{(1)} - x^{(0)})\|,$$

where the unavailable polynomial P^* minimizes $\|P(T)(x^{(1)} - x^{(0)})\|$ amongst all polynomials P of degree exactly $k - 1$ whose coefficients sum to 1.

The quality of the bound (in particular, its eventual convergence to 0) crucially depends on $\|P(T)\|$. Using the Crouzeix conjecture (Crouzeix, 2004) Bollapragada et al. (2018) managed to bound $\|P(T)\|$, with P a polynomial:

$$\|P(T)\| \leq c \max_{z \in W(T)} |P(z)|, \quad (10)$$

with $c \geq 2$ (Crouzeix, 2007; Crouzeix and Palencia, 2017), and $W(T)$ the numerical range:

$$W(T) \triangleq \{x^* T x : \|x\|_2 = 1, x \in \mathbb{C}^p\}. \quad (11)$$

Since there is no general formula for this bound, Bollapragada et al. (2018) used numerical bounds on

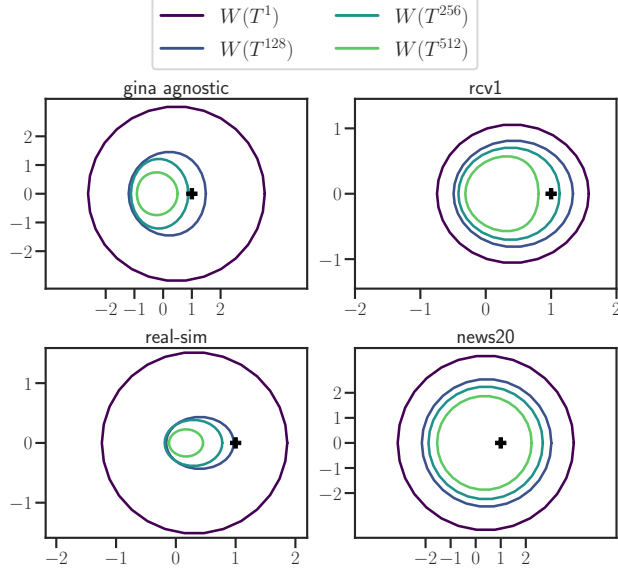


Figure 2: Numerical range of T^q as q varies; T is the iteration matrix of Ridge regression problems with conditioning $\kappa = 10^3$, on 4 datasets. The black cross marks the $(1, 0)$ point, which should lie outside the range for the theoretical bound to be useful.

$W(T^q)$ to ensure convergence. Figure 2 displays the numerical range $W(T^q)$ in the complex plane for $q \in \{1, 128, 256, 512\}$. In order to be able to apply the theoretical result from [Bollapragada et al. \(2018\)](#), one must chose q such that the point $(1, 0)$ is not contained in $W(T^q)$, and extrapolate $x^{(0)}, x^{(q)}, x^{(2q)}, \dots$. One can see on Figure 2 that large values of q are needed, unusable in practice: $q = 512$ is greater than the number of iterations needed to converge on some problems. Moreover, Anderson acceleration seems to provide speed up on coordinate descent even with $q = 1$ as we perform, which highlights the need for refined bounds for Anderson acceleration on nonsymmetric matrices.

We propose two means to fix this lack of theoretical results: to modify the algorithm in order to have a more amenable iteration matrix ([Section 2.4](#)), or to perform a simple cost function decrease check ([Section 2.5](#)).

2.4 Pseudo-symmetrization of T

A first idea to make coordinate descent theoretically amenable to extrapolation is to perform updates of coefficients from indices 1 to p , followed by a reversed pass from p to 1. This leads to an iteration matrix which is not symmetric either but friendlier: it writes

$$T^{\text{CD-sym}} \triangleq H^{-1/2} S H^{1/2}, \quad (12)$$

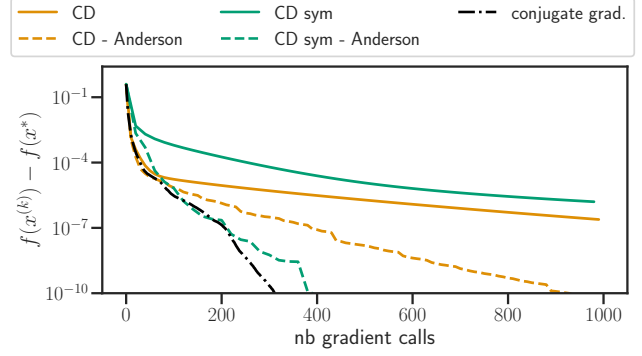


Figure 3: **OLS, rcv1**. Suboptimality as a function of the number of gradient calls on the 5000 first columns of the dataset *rcv1*.

with

$$S = \left(\text{Id}_p - H^{1/2} \frac{e_1 e_1^\top}{H_{11}} H^{1/2} \right) \times \dots \times \left(\text{Id}_p - H^{1/2} \frac{e_p e_p^\top}{H_{pp}} H^{1/2} \right) \\ \times \left(\text{Id}_p - H^{1/2} \frac{e_p e_p^\top}{H_{pp}} H^{1/2} \right) \times \dots \times \left(\text{Id}_p - H^{1/2} \frac{e_1 e_1^\top}{H_{11}} H^{1/2} \right). \quad (13)$$

S is symmetric, thus, S and T (which has the same eigenvalues as S), are diagonalisable with real eigenvalues. We call these iterations pseudo-symmetric, and show that this structure allows to preserve the guarantees of Anderson extrapolation.

Proposition 3 (Pseudosym. $T = H^{-1/2} S H^{1/2}$). Let T be the iteration matrix of pseudo-symmetric coordinate descent: $T = H^{-1/2} S H^{1/2}$, with S the symmetric semi-definite positive matrix of (12). Let x^* be the limit of the sequence $(x^{(k)})$. Let $\zeta = (1 - \sqrt{1 - \rho}) / (1 + \sqrt{1 - \rho})$. Then $\rho = \rho(T) = \rho(S) < 1$ and the iterates of offline Anderson acceleration satisfy:

$$\|x_{\text{e-off}}^{(k)} - x^*\|_B \leq \sqrt{\kappa(H)} \frac{2\zeta^{k-1}}{1 + \zeta^{2(k-1)}} \|x^{(0)} - x^*\|_B, \quad (14)$$

and thus those of online extrapolation satisfy:

$$\|x_{\text{e-on}}^{(k)} - x^*\|_B \leq \left(\sqrt{\kappa(H)} \frac{2\zeta^{K-1}}{1 + \zeta^{2(K-1)}} \right)^{k/K} \|x^{(0)} - x^*\|_B. \quad (15)$$

Proof of [Proposition 3](#) can be found in [Appendix B](#). [Proposition 3](#) shows accelerated convergence rates for the offline Anderson acceleration, but a $\sqrt{\kappa(H)}$ appears in the rate of the online Anderson acceleration, meaning that K must be large enough that ζ^K mitigates this effect. This factor however seems like a theoretical artefact of the proof, since we observed significant speed up of the online Anderson acceleration, even with bad conditioning of H (see [Figure 3](#)).

Algorithm 3 Online Anderson PCD (proposed)

```

init:  $x^{(0)} \in \mathbb{R}^p$ 
1 for  $k = 1, \dots$  do
2    $x = x^{(k-1)}$ 
3   for  $j = 1, \dots, p$  do
4      $\tilde{x}_j = x_j$ 
5      $x_j = \text{prox}_{\frac{\lambda}{L_j} g_j}(x_j - A_{:,j}^\top \nabla f(Ax)/L_j)$ 
6      $Ax += (x_j - \tilde{x}_j)A_{:,j}$ 
7    $x^{(k)} = x$  // regular iter.  $\mathcal{O}(np)$ 
8   if  $k = 0 \bmod K$  then // extrapol.,  $\mathcal{O}(K^3 + pK^2)$ 
9      $U = [x^{(k-K+1)} - x^{(k-K)}, \dots, x^{(k)} - x^{(k-1)}]$ 
10     $c = (U^\top U)^{-1} \mathbf{1}_K / \mathbf{1}_K^\top (U^\top U)^{-1} \mathbf{1}_K \in \mathbb{R}^K$ 
11     $x_e = \sum_{i=1}^K c_i x^{(k-K+i)}$ 
12    if  $f(Ax_e) + \lambda g(x_e) \leq f(x^{(k)}) + \lambda g(x^{(k)})$  then
13       $x^{(k)} = x_e$  // guaranteed convergence
14 return  $x^{(k)}$ 

```

Figure 3 illustrates the convergence speed of cyclic and pseudo-symmetric coordinate descent on the *rcv1* dataset. Anderson acceleration provides speed up for both versions. Interestingly, on this quadratic problem, the non extrapolated pseudo-symmetric iterations perform poorly, worse than cyclic coordinate descent. However, the performances are reversed for their extrapolated counterparts: the pseudo-symmetrized version is better than the cyclic one (which has a nonsymmetric iteration matrix). Finally, Anderson extrapolation on the pseudo-symmetrized version even reaches the conjugate gradient performance.

2.5 Generalization to nonquadratic and proposed algorithm

After devising and illustrating an Anderson extrapolated coordinate descent procedure for a simple quadratic objective, our goal is to apply Anderson acceleration on problems where coordinate descent achieve state-of-the-art results, *i.e.*, of the form:

$$\min_{x \in \mathbb{R}^p} f(Ax) + \lambda g(x) := f(Ax) + \lambda \sum_{j=1}^p g_j(x_j) \quad , \quad (16)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, γ -smooth and g_j 's are proper, closed and convex functions. As examples, we allow $g = 0$, $g = \|x\|_1$, $g = \frac{1}{2}\|x\|_2^2$, $g = \|x\|_1 + \frac{\rho}{2\lambda}\|x\|^2$. One pass of proximal coordinate descent from 1 to p can be seen as a nonlinear fixed point iteration:

$$x^{(k+1)} = \psi(x^{(k)}) \quad . \quad (17)$$

Proposition 4. If f is convex and smooth and \mathcal{C}^2 , g_j are convex smooth and \mathcal{C}^2 , then ψ is differentiable and

$$x^{(k+1)} = D\psi(x^*)(x^{(k)} - x^*) + x^* + o(\|x^{(k)} - x^*\|) \quad .$$

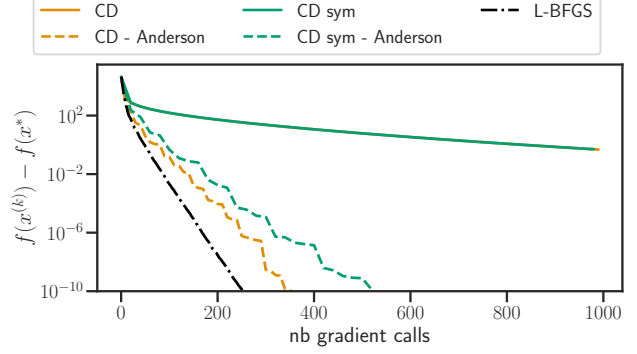


Figure 4: ℓ_2 -regularised logistic regression, *real-sim*. Suboptimality as a function of the number of gradient calls on the 2000 first features of the *real-sim* dataset, Tikhonov strength set so that $\kappa = 10^5$.

Therefore, iterations of proximal coordinate descent for this problem lead to noisy linear iterations. Proof of Proposition 4 can be found in Appendix B. Figure 4 shows the performance of Anderson extrapolation on a ℓ^2 -regularised logistic regression problem:

$$\arg \min_{x \in \mathbb{R}^p} \sum_{i=1}^n \log(1 + e^{-y_i A_{i,:} x}) + \frac{\lambda}{2} \|x\|_2^2 \quad . \quad (18)$$

One can see that despite the better theoretical properties of the pseudo-symmetrized coordinate descent, Anderson acceleration on coordinate descent seems to work better on the cyclic coordinate descent. We thus choose to apply Anderson extrapolation on the cyclic coordinate descent (Algorithm 3), while adding a step checking the decrease of the objective function in order to ensure convergence. Finally, we can also use Algorithm 3 in the non smooth case where $g = \|\cdot\|_1$, since coordinate descent achieves support identification when the solution is unique, after which the objective becomes differentiable. There is therefore a linear structure after a sufficient number of iterations (Massias et al., 2019, Prop. 10).

3 Experiments

An implementation relying on numpy, numba and cython (Harris et al., 2020; Lam et al., 2015; Behnel et al., 2011), with scripts to reproduce the figures, is available at <https://mathurinm.github.io/andersoncd>

We first show how we set the hyperparameters of Anderson extrapolation (Section 3.1). Then we show that Anderson extrapolation applied to proximal coordinate descent outperforms other first order algorithms on standard Machine Learning problems (Section 3.2).

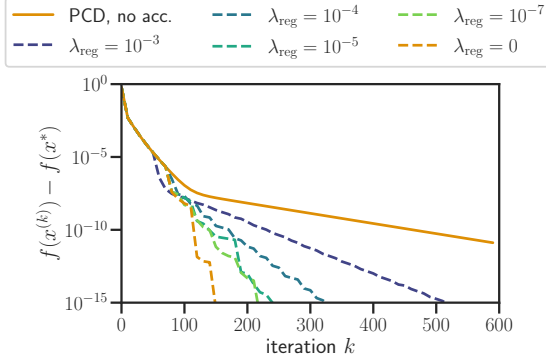


Figure 5: **Influence of λ_{reg} , sparse logistic regression, *rcv1*.** Influence of the regularization amount when solving a sparse logistic regression using Anderson extrapolation with proximal coordinate descent (PCD) on the *rcv1* dataset, $K = 5$, $\lambda = \lambda_{\text{max}}/30$.

3.1 Parameter setting

Anderson extrapolation relies on 2 hyperparameters: the number of extrapolated points K , and the amount of regularization eventually used when solving the linear system to obtain the coefficients $c \in \mathbb{R}^K$. Based on the conclusions of this section, we fix these parameters for all the subsequent experiments in Section 3.2: *no regularization and $K = 5$* .

Influence of the regularization. Scieur et al. (2016) provided accelerated complexity rates for *regularized* Anderson extrapolation: a term $\lambda_{\text{reg}} \|c\|^2$ is added to the objective of Equation (2). The closed-form formula for the coefficients is then $(U^\top U + \lambda_{\text{reg}} \text{Id}_K)^{-1} \mathbf{1}_K / \mathbf{1}_K^\top (U^\top U + \lambda_{\text{reg}} \text{Id}_K)^{-1} \mathbf{1}_K$.

However, similarly to Mai and Johansson (2019) and Poon and Liang (2020) we observed that regularizing the linear system does not seem necessary, and can even hurt the convergence speed. Figure 5 shows the influence of the regularization parameter on the convergence on the *rcv1* dataset for a sparse logistic regression problem, with $K = 5$ and $\lambda = \lambda_{\text{max}}/30$. The more the optimization problem is regularized, the more the convergence speed is deteriorated. We observed the same phenomenon when solving least squares problems (Figure 10, Appendix A.1). Thus *we choose not to regularize* when solving the linear system for the extrapolation coefficients. We simply check if the extrapolated point yields a lower objective function than the current regular iterate (see Algorithm 3).

Influence of K . Figure 6 shows the impact of K on the convergence speed. Although the performance depends on K , it seems that the dependency is loose, as for $K \in \{5, 10, 20\}$ the acceleration is roughly the

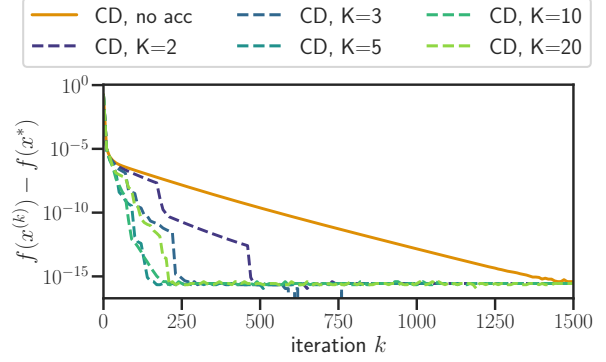


Figure 6: **Influence of K , quadratic, *rcv1*.** Influence of the number of iterates K used to perform Anderson extrapolation with coordinate descent (CD) on a quadratic with the *rcv1* dataset (1000 first columns).

same. Therefore, we do not treat K as a parameter and fix it to $K = 5$. Every K iterations Anderson accelerated algorithms require to solve a $K \times K$ linear system. For $K = 5$ it is marginal compared to a gradient call: *i.e.*, $5^3 + 5^2 p \ll np$ in our settings.

3.2 Numerical comparison on Machine Learning problems

We compare multiple algorithms to solve popular Machine Learning problems: the Lasso, the elastic net, and sparse logistic regression (experiments on group lasso are in Appendix A.5). The compared algorithms are the following: proximal gradient descent (PGD, Combettes and Wajs 2005), Nesterov-like inertial PGD (FISTA, Beck and Teboulle 2009), Anderson accelerated PGD (Mai and Johansson, 2019; Poon and Liang, 2020), proximal coordinate descent (PCD, Tseng and Yun 2009), inertial PCD (Lin et al., 2014; Fercoq and Richtárik, 2015), Anderson accelerated PCD (ours). We use datasets from libsvm (Fan et al., 2008) and openml (Feurer et al., 2019) (Table 1), varying as much as possible to demonstrate the versatility of our approach. We also vary the convergence metric: we use suboptimality in the main paper, while graphs measuring the duality gaps are in Appendix A.

Lasso. Figure 7 shows the suboptimality $f(x^{(k)}) - f(x^*)$ of the algorithms on the Lasso problem:

$$x^* = \arg \min_{x \in \mathbb{R}^p} \frac{1}{2} \|y - Ax\|^2 + \lambda \|x\|_1, \quad (19)$$

as a function of the number of iterations for multiple datasets and values of λ . We parametrize λ as a fraction of $\lambda_{\text{max}} = \|A^\top y\|_\infty$, smallest regularization strength for which $x^* = 0$. Figure 7 highlights the

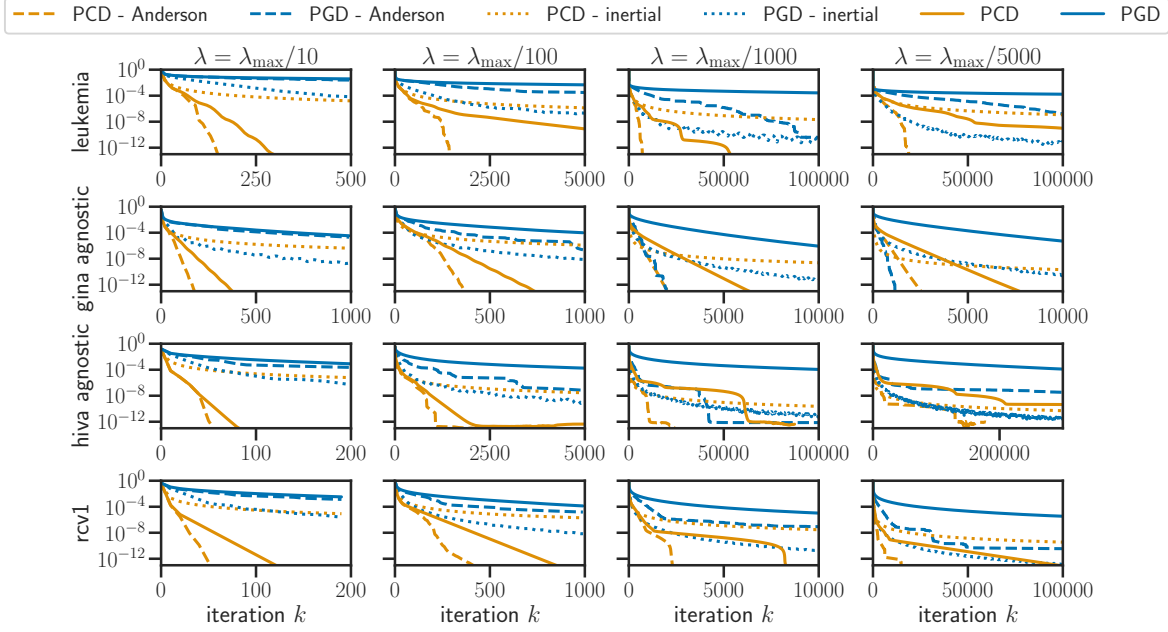


Figure 7: **Lasso, suboptimality.** Suboptimality $f(x^{(k)}) - f(x^*)$ as a function of the number of iterations for the Lasso on multiple datasets and values of λ .

Table 1: Datasets characteristics

name	n	p	density
<i>gina agnostic</i>	3468	970	1
<i>hiva agnostic</i>	4229	1617	1
<i>leukemia</i>	72	7129	1
<i>rcv1_train</i>	20 242	19 960	$3.7 \cdot 10^{-3}$
<i>real-sim</i>	72 309	20 958	$2.4 \cdot 10^{-3}$
<i>news20</i>	19 996	632 983	$6.1 \cdot 10^{-4}$

superiority of proximal coordinate descent over proximal gradient descent for Lasso problems on real-world datasets, and the benefits of extrapolation for coordinate descent. It shows that Anderson extrapolation can lead to a significant gain of performance. In particular Figure 7 shows that without restart, inertial coordinate descent (Lin et al., 2014; Fercoq and Richtárik, 2015) can slow down the convergence, despite its accelerated rate. Note that the smaller the value of λ , the harder the optimization: when λ decreases, more iterations are needed to reach a fixed suboptimality. The smaller λ is (*i.e.*, the harder the problem), the more efficient Anderson extrapolation is.

As in Mai and Johansson (2019), Anderson PGD performs well when $p < n$: Anderson PGD outperforms Anderson PCD on the *gina agnostic* dataset, where $p < n$. On all the other datasets, especially when $p > n$ and when the values of λ are large, An-

derson PCD outperforms Anderson PGD.

Other convergence metrics can be considered, since $f(x^*)$ is unknown to the practitioner: for the Lasso, it is also common to use the duality gap as a stopping criterion (Massias et al., 2018). Thus, for completeness, we provide Figure 11 in appendix, which shows the duality gap as a function of the number of iterations. With this metric of convergence, Anderson PCD also significantly outperforms its competitors.

Elastic net. Anderson extrapolation is easy to extend to other estimators than the Lasso. Figure 8 (and Figure 12 in appendix) show the superiority of the Anderson extrapolation approach over proximal gradient descent and its accelerated version for the elastic net problem (Zou and Hastie, 2005):

$$\arg \min_{x \in \mathbb{R}^p} \frac{1}{2n} \|y - Ax\|^2 + \lambda \|x\|_1 + \frac{\rho}{2} \|x\|_2^2. \quad (20)$$

In particular, we observe that the more difficult the problem, the more useful the Anderson extrapolation: it is visible on Figures 8 and 12 that going from $\rho = \lambda/10$ to $\rho = \lambda/100$ lead to an increase in the number of iterations to achieve similar suboptimality for the classical coordinate descent, whereas the impact is more limited on the coordinate descent with Anderson extrapolation.

Finally, for a nonquadratic data-fit, here sparse logistic regression, we still demonstrate the applicability of extrapolated coordinate descent.

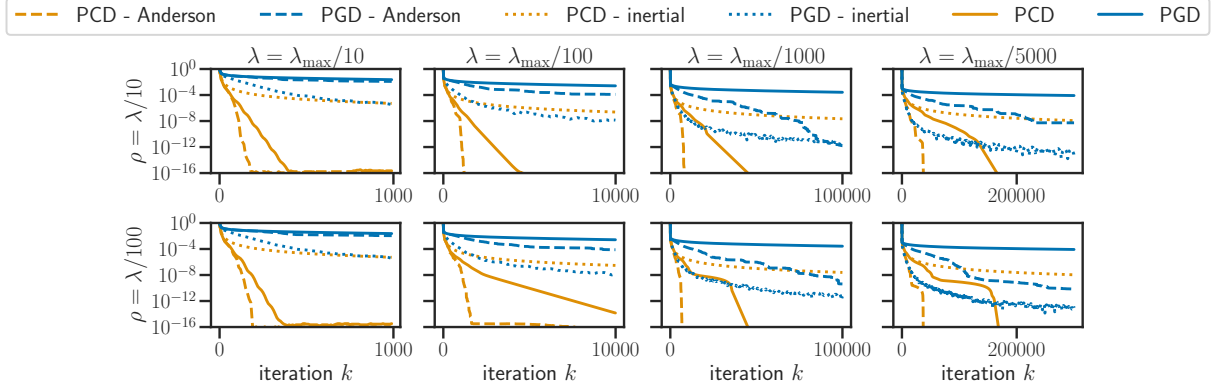


Figure 8: **Enet, suboptimality.** Suboptimality $f(x^{(k)}) - f(x^*)$ as a function of the number of iterations for the elastic net on Leukemia dataset, for multiple values of λ and ρ .

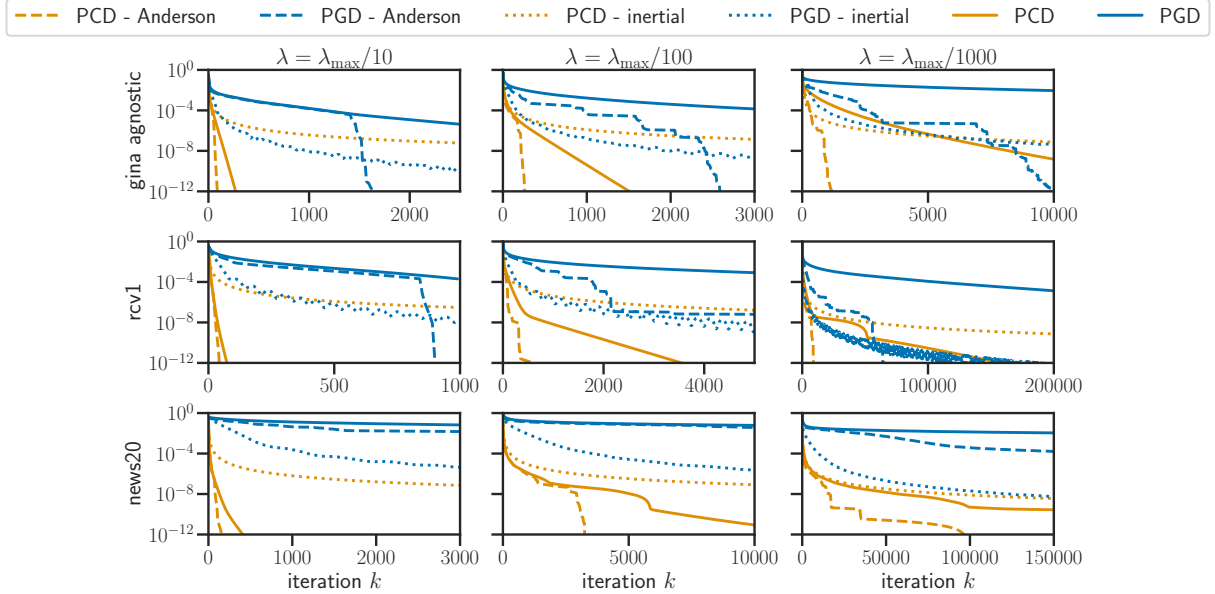


Figure 9: **ℓ_1 -regularised logistic regression, suboptimality.** Suboptimality $f(x^{(k)}) - f(x^*)$ as a function of the number of iterations for ℓ_1 -regularized logistic regression on multiple datasets and values of λ .

Sparse logistic regression. Figure 9 represents the suboptimality as a function of the number of iterations on a sparse logistic regression problem:

$$\arg \min_{x \in \mathbb{R}^p} \sum_{i=1}^n \log(1 + e^{-y_i A_{i,:} x}) + \lambda \|x\|_1, \quad (21)$$

for multiple datasets and values of λ . We parametrize λ as a fraction of $\lambda_{\max} = \|A^\top y\|_\infty / 2$. As for the Lasso and the elastic net, the smaller the value of λ , the harder the problem and Anderson CD outperforms its competitors.

Conclusion In this work, we have proposed to accelerate coordinate descent using Anderson extrapolation. We have exploited the fixed point itera-

tions followed by coordinate descent iterates on multiple Machine Learning problems to improve their convergence speed. We have circumvented the non-symmetry of the iteration matrices by proposing a pseudo-symmetric version for which accelerated convergence rates have been derived. In practice, we have performed an extensive validation to demonstrate large benefits on multiple datasets and problems of interests. For future works, the excellent performance of Anderson extrapolation for cyclic coordinate descent calls for a more refined analysis of the known bounds, through a better analysis of the spectrum and numerical range of the iteration matrices.

References

- D. G. Anderson. Iterative procedures for nonlinear integral equations. *Journal of the ACM*, 12(4):547–560, 1965.
- M. Barré, A. Taylor, and A. d’Aspremont. Convergence of constrained anderson acceleration. *arXiv preprint arXiv:2010.15482*, 2020.
- A. Beck. *First-Order Methods in Optimization*, volume 25. SIAM, 2017.
- A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202, 2009.
- S. Behnel, R. Bradshaw, C. Citro, L. Dalcin, D. S. Seljebotn, and K. Smith. Cython: The best of both worlds. *Computing in Science Engineering*, 13(2):31–39, 2011.
- Q. Bertrand, Q. Klopfenstein, M. Blondel, S. Vaiter, A. Gramfort, and J. Salmon. Implicit differentiation of lasso-type models for hyperparameter optimization. *ICML*, 2020.
- M. Blondel and F. Pedregosa. Lightning: large-scale linear classification, regression and ranking in python, 2016.
- R. Bollapragada, D. Scieur, and A. d’Aspremont. Non-linear acceleration of momentum and primal-dual algorithms. *arXiv preprint arXiv:1810.04539*, 2018.
- C. Brezinski, M. Redivo-Zaglia, and Y. Saad. Shanks sequence transformations and anderson acceleration. *SIAM Review*, 60(3):646–669, 2018.
- P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Modeling & Simulation*, 4(4):1168–1200, 2005.
- M. Crouzeix. Bounds for analytical functions of matrices. *Integral Equations and Operator Theory*, 48(4):461–477, 2004.
- M. Crouzeix. Numerical range and functional calculus in hilbert space. *Journal of Functional Analysis*, 244(2):668–690, 2007.
- M. Crouzeix and C. Palencia. The numerical range is a $(1+2)$ -spectral set. *SIAM Journal on Matrix Analysis and Applications*, 38(2):649–655, 2017.
- R. P. Eddy. Extrapolating to the limit of a vector sequence. In *Information linkage between applied mathematics and industry*, pages 387–396. Elsevier, 1979.
- V. Eyert. A comparative study on methods for convergence acceleration of iterative vector sequences. *Journal of Computational Physics*, 124(2):271–285, 1996.
- R. E. Fan, K. W. Chang, C. J. Hsieh, X. R. Wang, and C. J. Lin. Liblinear: A library for large linear classification. *JMLR*, 9:1871–1874, 2008.
- O. Fercoq and P. Richtárik. Accelerated, parallel, and proximal coordinate descent. *SIAM Journal on Optimization*, 25(4):1997–2023, 2015.
- M. Feurer, J. N. van Rijn, A. Kadra, P. Gijsbers, N. Mallik, S. Ravi, A. Müller, J. Vanschoren, and F. Hutter. Openml-python: an extensible python api for openml. *arXiv:1911.02490*, 2019.
- J. Friedman, T. Hastie, and R. Tibshirani. glmnet: Lasso and elastic-net regularized generalized linear models. *R package version*, 1(4), 2009.
- J. Friedman, T. J. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.*, 33(1):1, 2010.
- A. Fu, J. Zhang, and S. Boyd. Anderson accelerated Douglas-Rachford splitting. *arXiv preprint arXiv:1908.11482*, 2019.
- G. H. Golub and R. S. Varga. Chebyshev semi-iterative methods, successive overrelaxation iterative methods, and second order richardson iterative methods. *Numerische Mathematik*, 3(1):147–156, 1961.
- R. Gribonval and M. Nikolova. A characterization of proximity operators. *Journal of Mathematical Imaging and Vision*, 62(6):773–789, 2020.
- C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, et al. Array programming with numpy. *arXiv preprint arXiv:2006.10256*, 2020.
- M. R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49(6):409–436, 1952.
- Quentin Klopfenstein, Quentin Bertrand, Alexandre Gramfort, Joseph Salmon, and Samuel Vaiter. Model identification and local linear convergence of coordinate descent. *arXiv preprint arXiv:2010.11825*, 2020.
- S. K. Lam, A. Pitrou, and S. Seibert. Numba: A llvm-based python jit compiler. In *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*, pages 1–6, 2015.

- Q. Lin, Z. Lu, and L. Xiao. An accelerated proximal coordinate gradient method. In *NeurIPS*, pages 3059–3067. 2014.
- V. V. Mai and M. Johansson. Anderson acceleration of proximal gradient methods. In *ICML*. 2019.
- J. Mairal. *Sparse coding for machine learning, image processing and computer vision*. PhD thesis, École normale supérieure de Cachan, 2010.
- M. Massias, A. Gramfort, and J. Salmon. Celer: a fast solver for the lasso with dual extrapolation. 2018.
- M. Massias, S. Vaiter, A. Gramfort, and J. Salmon. Dual extrapolation for sparse generalized linear models. *arXiv preprint arXiv:1907.05830*, 2019.
- Y. Nesterov. A method for solving a convex programming problem with rate of convergence $O(1/k^2)$. *Soviet Math. Doklady*, 269(3):543–547, 1983.
- Y. Nesterov. *Introductory lectures on convex optimization*, volume 87 of *Applied Optimization*. Kluwer Academic publishers, Boston, MA, 2004.
- Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- A. Y. Ng. Feature selection, l_1 vs. l_2 regularization, and rotational invariance. In *ICML*, page 78, 2004.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *JMLR*, 12:2825–2830, 2011.
- C. Poon and J. Liang. Trajectory of alternating direction method of multipliers and adaptive acceleration. In *NeurIPS*, pages 7357–7365, 2019.
- C. Poon and J. Liang. Geometry of first-order methods and adaptive acceleration. *arXiv preprint arXiv:2003.03910*, 2020.
- P. Pulay. Convergence acceleration of iterative sequences. the case of scf iteration. *Chemical Physics Letters*, 73(2):393–398, 1980.
- D. Scieur. Generalized framework for nonlinear acceleration. *arXiv preprint arXiv:1903.08764*, 2019.
- D. Scieur, A. d’Aspremont, and F. Bach. Regularized nonlinear acceleration. In *Advances In Neural Information Processing Systems*, pages 712–720, 2016.
- S. Shalev-Shwartz and T. Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. *arXiv preprint arXiv:1309.2375*, 2013.
- H.-J. Shi, S. Tu, Y. Xu, and W. Yin. A primer on coordinate descent algorithms. *arXiv preprint arXiv:1610.00040*, 2016.
- A. Sidi. *Vector extrapolation methods with applications*. SIAM, 2017.
- D. A. Smith, W. F. Ford, and A. Sidi. Extrapolation methods for vector sequences. *SIAM review*, 29(2):199–233, 1987.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 58(1):267–288, 1996.
- P. Tseng and S. Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1-2):387–423, 2009.
- S. J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.
- P. Wynn. Acceleration techniques for iterated vector and matrix problems. *Mathematics of Computation*, 16(79):301–322, 1962.
- J. Zhang, B. O’Donoghue, and S. Boyd. Globally convergent type-I Anderson acceleration for non-smooth fixed-point iterations. *arXiv preprint arXiv:1808.03971*, 2018.
- H. Zou and T. J. Hastie. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 67(2):301–320, 2005.

A Additional experiments

In this section, we include the counterparts of Figures 7 to 9, but display the duality gap instead of the suboptimality. Indeed, since x^* is not available in practice, the suboptimality cannot be used as a stopping criterion. To create a dual feasible point, we use the classical technique of residual rescaling (Mairal, 2010).

A.1 OLS

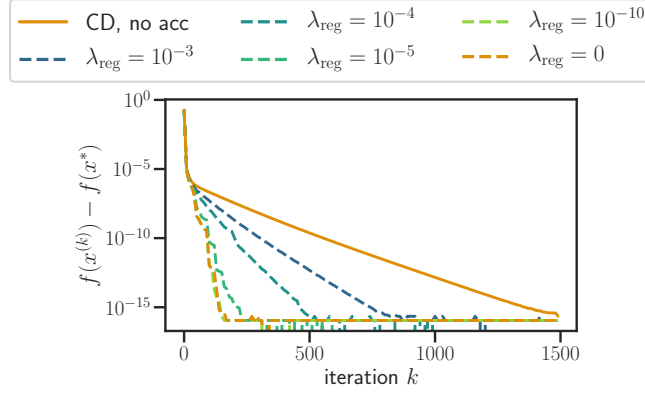


Figure 10: **Influence of λ_{reg} , quadratic, *rcv1*.** Influence of the regularization amount when solving the linear system for Anderson extrapolation with coordinate descent (CD) on a quadratic with the *rcv1* dataset (1000 first columns), $K = 5$.

A.2 Lasso

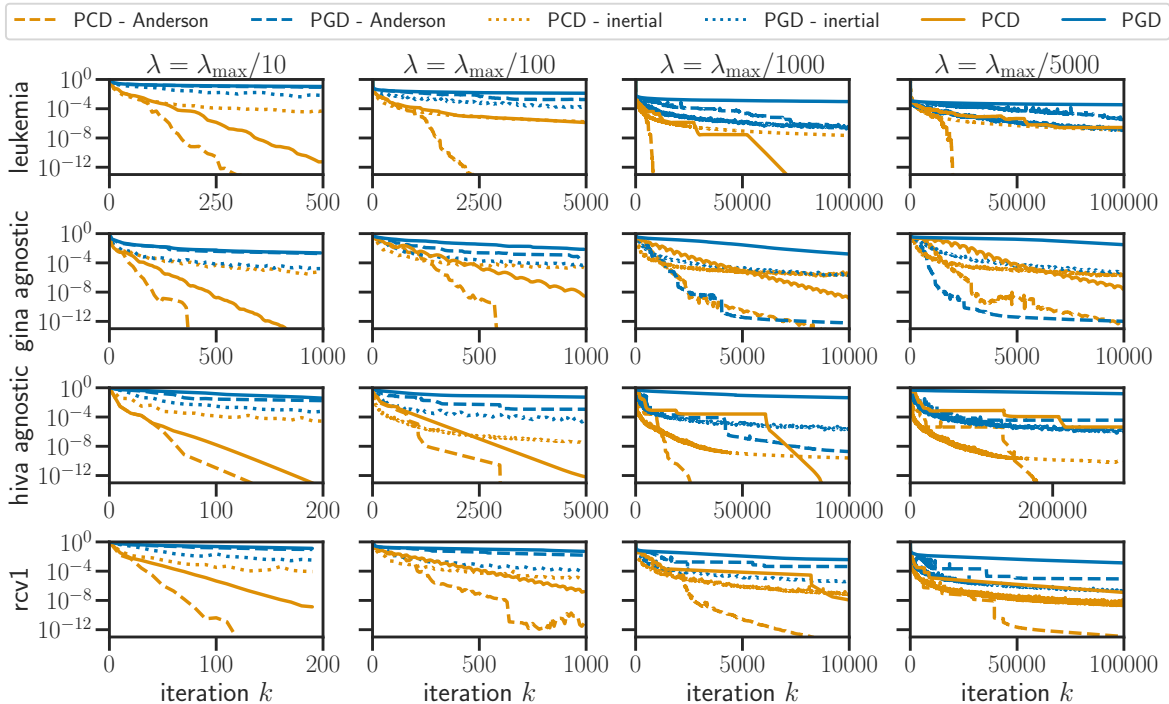


Figure 11: **Lasso, duality gap.** Duality gap along iterations for the Lasso on various datasets and values of λ .

A.3 Elastic net

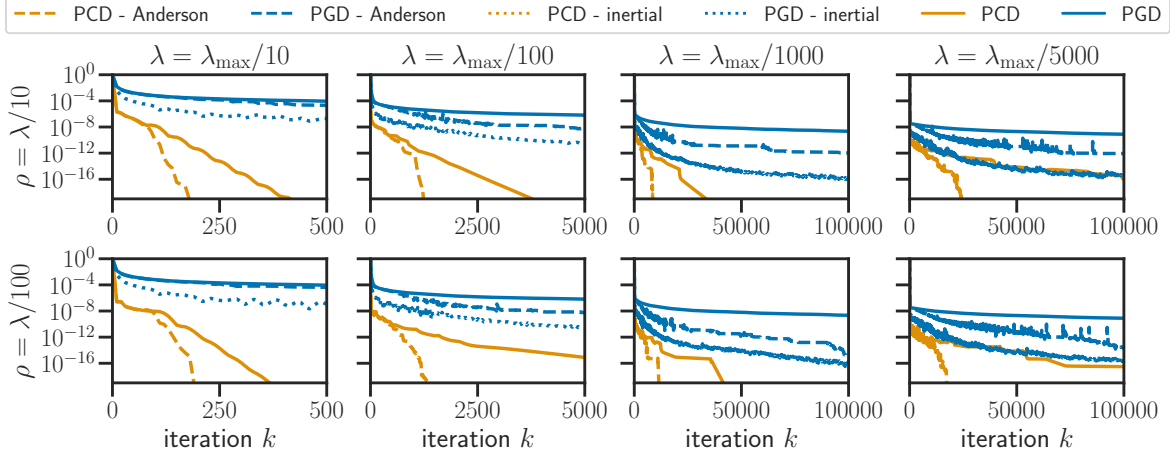


Figure 12: **Elastic net, duality gap.** Duality gap as a function of the number of iterations for the elastic net on Leukemia dataset, for multiple values of λ and ρ .

A.4 Sparse logistic regression

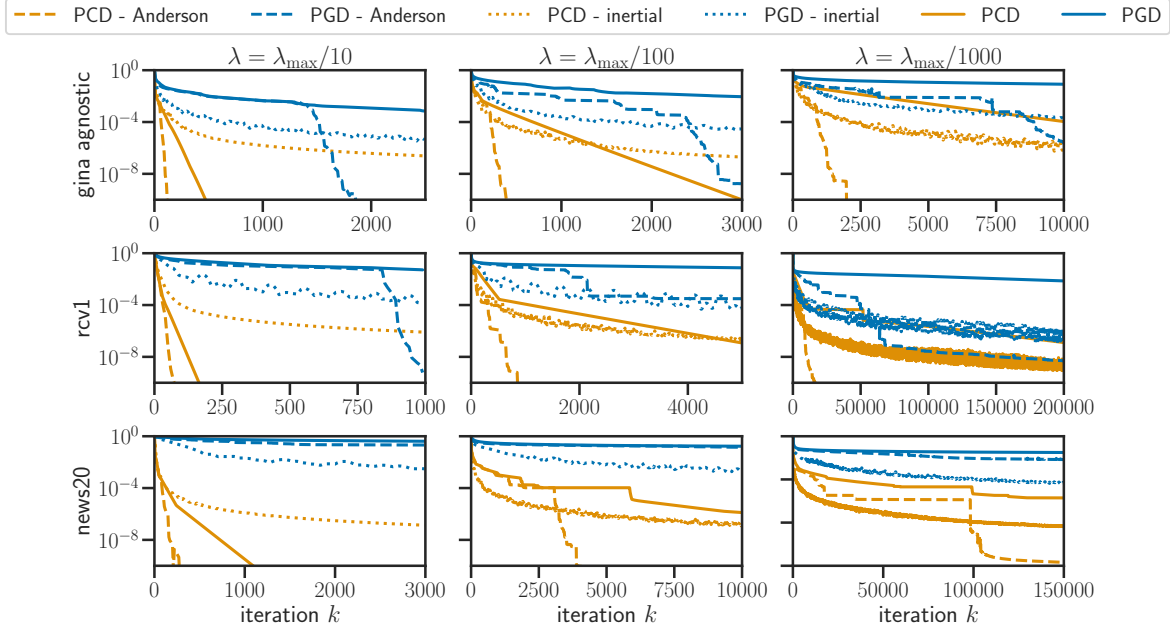


Figure 13: **ℓ_1 -regularised logistic regression, duality gap.** Duality gap as a function of the number of iterations for ℓ_1 -regularized logistic regression on multiple datasets and values of λ .

A.5 Group Lasso

In this section we consider the group Lasso, with a design matrix $A \in \mathbb{R}^{n \times p}$, a target $y \in \mathbb{R}^n$, and a partition \mathcal{G} of $[p]$ (elements of the partition being the disjoint groups):

$$\arg \min_{x \in \mathbb{R}^p} \frac{1}{2} \|y - Ax\|^2 + \lambda \sum_{g \in \mathcal{G}} \|x_g\|, \quad (22)$$

where for $g \in \mathcal{G}$, $x_g \in \mathbb{R}^{|g|}$ is the subvector of x composed of coordinates in g . the group Lasso can be solved via Proximal Gradient Descent and by Block Coordinate Descent (BCD), the latter being amenable to Anderson Acceleration. As Figure 14 shows, the superiority of Anderson accelerated block coordinate descent is on par with the one observed on the problems studied above.

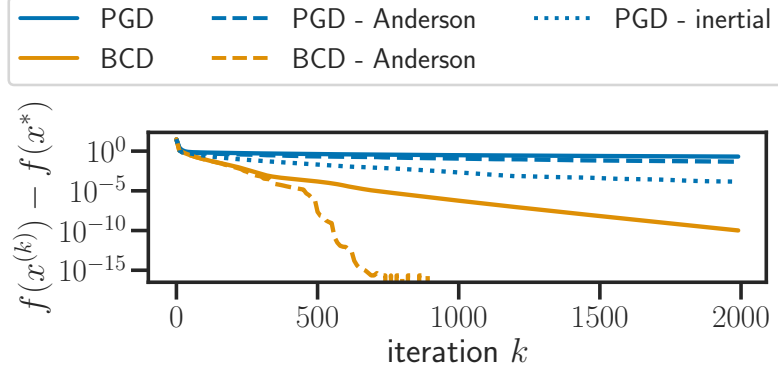


Figure 14: **Group Lasso, suboptimality.** Suboptimality as a function of the number of iterations for the group Lasso on the *Leukemia* dataset, $\lambda = \lambda_{\max}/100$. Groups are artificially taken as consecutive blocks of 5 features.

B Proofs of Propositions 3 and 4

B.1 Proofs of Proposition 3

Lemma 5. First we link the quantity computed in Equation (2) to the extrapolated quantity $\sum_{i=1}^k c_i x^{(i-1)}$. For all $c \in \mathbb{R}^k$ such that $\sum_{i=1}^k c_i = 1$:

$$\sum_{i=1}^k c_i (x^{(i)} - x^{(i-1)}) = (T - \text{Id}) \left(\sum_{i=1}^k c_i x^{(i-1)} - x^* \right) . \quad (23)$$

Proof. Since $x^{(i)} = Tx^{(i-1)} + (x^* - Tx^*)$,

$$\begin{aligned} c_i (x^{(i)} - x^{(i-1)}) &= c_i (Tx^{(i-1)} + x^* - Tx^* - x^{(i-1)}) \\ &= (T - \text{Id}) c_i (x^{(i-1)} - x^*) . \end{aligned} \quad (24)$$

Hence, since $\sum_{i=1}^k c_i = 1$,

$$\sum_{i=1}^k c_i (x^{(i)} - x^{(i-1)}) = (T - \text{Id}) \left(\sum_{i=1}^k c_i x^{(i-1)} - x^* \right) . \quad (25)$$

□

Lemma 6. For all $c \in \mathbb{R}^k$ such that $\sum_{i=1}^k c_i = 1$,

$$\|(T - \text{Id})(x_{\text{e-off}}^{(k)} - x^*)\| \leq \sqrt{\kappa(H)} \left\| \sum_{i=0}^{k-1} c_i S^i \right\| \|(T - \text{Id})(x^{(0)} - x^*)\| . \quad (26)$$

Proof. In this proof, we denote by c^* the solution of (2). We use the fact that for all $c \in \mathbb{R}^k$ such that $\sum_{i=1}^k c_i = 1$,

$$\left\| \sum_{i=1}^k c_i^* (x^{(i)} - x^{(i-1)}) \right\| = \min_{\substack{c \in \mathbb{R}^k \\ \sum_i c_i = 1}} \left\| \sum_{i=1}^k c_i (x^{(i)} - x^{(i-1)}) \right\| \leq \left\| \sum_{i=1}^k c_i (x^{(i)} - x^{(i-1)}) \right\| . \quad (27)$$

Then we use twice [Lemma 5](#) for the left-hand and right-hand side of [Equation \(27\)](#). Using [Lemma 5](#) with the c_i^* minimizing [Equation \(2\)](#) we have for all $c_i \in \mathbb{R}$ such that $\sum_{i=1}^k c_i = 1$:

$$\begin{aligned}
\|(T - \text{Id})(x_e - x^*)\| &= \left\| \sum_{i=1}^k c_i^* (x^{(i)} - x^{(i-1)}) \right\| \\
&\leq \left\| \sum_{i=1}^k c_i (x^{(i)} - x^{(i-1)}) \right\| \\
&= \|(T - \text{Id}) \sum_{i=1}^k c_i (x^{(i-1)} - x^{(*)})\| \\
&= \|(T - \text{Id}) \sum_{i=1}^k c_i T^{i-1} (x^{(0)} - x^*)\| \\
&= \left\| \sum_{i=1}^k c_i T^{i-1} (T - \text{Id})(x^{(0)} - x^*) \right\| \\
&= \left\| \sum_{i=1}^k c_i T^{i-1} \right\| \times \|(T - \text{Id})(x^{(0)} - x^*)\| \\
&\leq \|H^{-1/2} \sum_{i=1}^k c_i S^{i-1} H^{1/2}\| \times \|(T - \text{Id})(x^{(0)} - x^*)\| \\
&\leq \sqrt{\kappa(H)} \left\| \sum_{i=1}^k c_i S^{i-1} \right\| \times \|(T - \text{Id})(x^{(0)} - x^*)\| .
\end{aligned} \tag{28}$$

□

Proof. We apply [Lemma 6](#) by choosing c_i equal to the Chebyshev weights c_i^{Cb} . Using the proof of [Barré et al. \(2020, Prop. B. 2\)](#), we have, with $\zeta = \frac{1-\sqrt{1-\rho(T)}}{1+\sqrt{1-\rho(T)}}$:

$$\left\| \sum_{i=1}^k c_i^{\text{Cb}} S^{i-1} \right\| \leq \frac{2\zeta^{k-1}}{1+\zeta^{2(k-1)}} . \tag{29}$$

Combined with [Lemma 6](#) this concludes the proof:

$$\|(T - \text{Id})(x_e - x^*)\| \leq \sqrt{\kappa(H)} \left\| \sum_{i=1}^k c_i S^{i-1} \right\| \|(T - \text{Id})(x^{(0)} - x^*)\| \tag{30}$$

$$\leq \sqrt{\kappa(H)} \frac{2\zeta^{k-1}}{1+\zeta^{2(k-1)}} \|(T - \text{Id})(x^{(0)} - x^*)\| . \tag{31}$$

□

B.2 Proof of Proposition 4

Since g_j are \mathcal{C}^2 then prox_{g_j} are \mathcal{C}^1 , see [Gribonval and Nikolova \(2020, Cor. 1.b\)](#). Moreover, f is \mathcal{C}^2 and following [Massias et al. \(2019\)](#); [Klopfenstein et al. \(2020\)](#) we have that:

$$\psi_j : \mathbb{R}^p \rightarrow \mathbb{R}^p$$

$$x \mapsto \text{prox}_{g_j} \begin{pmatrix} x_1 \\ \vdots \\ x_{j-1} \\ \text{prox}_{\lambda g_j / L_j} (x_j - \gamma_j \nabla_j f(x)) \\ x_{j+1} \\ \vdots \\ x_p \end{pmatrix}, \quad (32)$$

is differentiable. Thus we have that the fixed point operator of coordinate descent: $\psi = \psi_p \circ \dots \circ \psi_1$ is differentiable. [Proposition 4](#) follows from the Taylor expansion of ψ in x^* .